# Recursive estimation in Riemannian manifolds

Salem Said 2019 - École d'été de Peyresq

CNRS – Université de Bordeaux

You should become able to read this :

Zhou & Said (2018) : fast asymptotically-efficient recursive estimation in a Riemannian manifold (https://arxiv.org/abs/1805.06811)

Nothing is more practical than a good theory !!

# Plan

Introduction with surfaces

### 2 Higher dimension

Convex stochastic optimisation

A Riemannian recursive estimation

### Proposed reading

### What is a manifold ?

Riemann (1854) : any space which one can describe by varying *n* real parameters {attitude of a solid body in space} {shape of a deformable elastic body} {color}

Abstract manifolds (1920s-30s) : the cartographer's definition

 $\{(z, w) \in \mathbb{C} \times \mathbb{C} | w^2 - z^2 + 4 = 0\}$  is a cylindre !!

a manifold is a space with at atlas which is a set of compatible local charts

♦ Whitney's theorem (1944) : all manifolds are concrete
 –any smooth *n*-dimensional manifold can be embedded into ℝ<sup>2n</sup>
 – the embedding is difficult, but we know it always exists



# What is a Riemannian manifold ?

Riemannian metric : length is measured by a quadratic form

think of an embedded manifold  $M \subset \mathbb{R}^N$ 

$$L(\gamma) = \int_0^1 \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle^{\frac{1}{2}} dt \quad \text{curve } \gamma : [0.1] \to M$$

 $\rightsquigarrow$  reparameterisation-invariant definition of length

Riemannian distance :

$$d(x, y) = \inf \{ L(\gamma); \gamma(0) = x \text{ and } \gamma(1) = y \}$$

♦ Geodesics :

locally length-minimising curves

 $L(\gamma \circ \phi) = L(\gamma)$ 

$$L(\gamma|[t, t + \epsilon]) = d(\gamma(t), \gamma(t + \epsilon)) \quad \text{for each } t \in (0, 1)$$

a geodesic may fail to be globally minimising !!

there exists an infinite choice of Riemannian metrics on a given manifold

Curvature is a more relevant quantity

Gauss-Bonnet for surfaces : 
$$\frac{1}{4\pi} \int_{M} R d\text{vol} = 2 - 2g$$



# What is a Riemannian manifold ?

Riemannian metric : length is measured by a quadratic form

#### think of an embedded manifold $M \subset \mathbb{R}^N$

$$L(\gamma) = \int_0^1 \langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle^{\frac{1}{2}} dt \quad \text{curve } \gamma : [0.1] \to M$$

 $\rightsquigarrow$  reparameterisation-invariant definition of length

Riemannian distance :

$$d(x, y) = \inf \{ L(\gamma); \gamma(0) = x \text{ and } \gamma(1) = y \}$$

**Oreconstantial Geodesics** :

locally length-minimising curves

 $L(\gamma \circ \phi) = L(\gamma)$ 

 $\begin{array}{l} \epsilon \text{ essentially depends} \\ on \ current \ point \ \gamma(t) \end{array} \qquad L(\gamma[[t, t+\epsilon]]) = d(\gamma(t), \gamma(t+\epsilon)) \qquad \text{for each } t \in (0, 1) \end{array}$ 

a geodesic may fail to be globally minimising !!

there exists an infinite choice of Riemannian metrics on a given manifold

Curvature is a more relevant quantity

Gauss-Bonnet for surfaces : 
$$\frac{1}{4\pi} \int_M R d\text{vol} = 2 - 2g$$



#### ♦ Geodesic spherical coordinates :

$$(r, \theta) \longmapsto \gamma_{\theta}(r)$$
 for  $r \in (0, \epsilon_x)$ ,  $\theta \in S^1$ 

♦ The metric :

$$\underbrace{\langle v, v \rangle = v_r^2 + f^2(r, \theta) v_{\theta}^2}_{\bullet}$$

$$\underbrace{ds^2 = dr^2 + f^2(r,\,\theta)\,d\theta^2}_{\underbrace{}}$$

scalar product

length element

♦ Jacobi equation : second order linear ode

 $f_{rr} + Kf = 0$   $K(r, \theta)$  sectional curvature

♦ Constant curvature :

$$\underbrace{f(r, \theta) = k^{-1} \sin(kr)}_{\text{positive curvature } k^2} \qquad \underbrace{f(r, \theta) = k^{-1} \sinh(kr)}_{\text{negative curvature } -k^2}$$

♦ Distance function (locally!!) :

$$y = \gamma_{\theta}(r) \implies d(x, y) = r$$

Hessian of distance function :

$$\nabla^2 r \simeq \left( \begin{array}{cc} 0 & 0 \\ 0 & f_r/f \end{array} \right)$$

#### ♦ Geodesic spherical coordinates :

 $\gamma_{\theta}$  unit speed geodesic in direction  $\theta$ 

 $(r, \theta) \longmapsto \gamma_{\theta}(r)$  for  $r \in (0, \epsilon_x)$ ,  $\theta \in S^1$ 

♦ The metric :

$$\underbrace{\langle v, v \rangle = v_r^2 + f^2(r, \theta) v_{\theta}^2}_{\bullet}$$

$$ds^2 = dr^2 + f^2(r,\,\theta)\,d\theta^2$$

scalar product

length element

♦ Jacobi equation : second order linear ode

 $f_{rr} + Kf = 0$   $K(r, \theta)$  sectional curvature

Constant curvature :

$$\underbrace{f(r, \theta) = k^{-1} \sin(kr)}_{\text{positive curvature } k^2} \qquad \underbrace{f(r, \theta) = k^{-1} \sinh(kr)}_{\text{negative curvature } -k^2}$$

♦ Distance function (locally!!) :

$$y = \gamma_{\theta}(r) \implies d(x, y) = r$$

♦ Hessian of distance function :

$$\nabla^2 r \simeq \left( \begin{array}{cc} 0 & 0 \\ 0 & f_r/f \end{array} \right)$$

#### ♦ Geodesic spherical coordinates :

 $\gamma_{\theta}$  unit speed geodesic in direction  $\theta$ 

 $(r, \theta) \mapsto \gamma_{\theta}(r)$ 

for  $r \in (0, \epsilon_x)$ ,  $\theta \in S^1$ 

♦ The metric :

$$\underbrace{\langle v, v \rangle = v_r^2 + f^2(r, \theta) v_{\theta}^2}_{\underbrace{}}$$

$$ds^2 = dr^2 + f^2(r,\,\theta)\,d\theta^2$$

scalar product

length element

♦ Jacobi equation : second order linear ode

K determines M up to coverings  $f_{rr} + Kf = 0$   $K(r, \theta)$  sectional curvature

Constant curvature :

$$\underbrace{f(r, \theta) = k^{-1} \sin(kr)}_{\text{positive curvature } k^2} \underbrace{f(r, \theta) = k^{-1} \sinh(kr)}_{\text{negative curvature } -k^2}$$

♦ Distance function (locally!!) :

$$y = \gamma_{\theta}(r) \implies d(x, y) = r$$

Hessian of distance function :

$$\nabla^2 r \simeq \left( \begin{array}{cc} 0 & 0 \\ 0 & f_r/f \end{array} \right)$$

4/29

#### Geodesic spherical coordinates :

 $\gamma_{\theta}$  unit speed geodesic in direction  $\theta$ 

 $(r, \theta) \mapsto \gamma_{\theta}(r)$ 

for  $r \in (0, \epsilon_x)$ ,  $\theta \in S^1$ 

♦ The metric :

$$\underbrace{\langle v, v \rangle = v_r^2 + f^2(r, \theta) v_{\theta}^2}_{\underbrace{}}$$

scalar product

length element

♦ Jacobi equation : second order linear ode

K determines M up to coverings  $f_{rr} + Kf = 0$   $K(r, \theta)$  sectional curvature

Constant curvature :

 $f(r, \theta) \sim r \text{ for } \underbrace{f(r, \theta) = k^{-1} \sin(kr)}_{\text{positive curvature } k^2} \qquad \underbrace{f(r, \theta) = k^{-1} \sinh(kr)}_{\text{negative curvature } -k^2}$ 

◇ Distance function (locally!!) :

$$y = \gamma_{\theta}(r) \implies d(x, y) = r$$

Hessian of distance function :

$$\nabla^2 r \simeq \left( \begin{array}{cc} 0 & 0 \\ 0 & f_r/f \end{array} \right)$$

#### ♦ Geodesic spherical coordinates :

 $\gamma_{\theta}$  unit speed geodesic in direction  $\theta$ 

 $(r, \theta) \mapsto \gamma_{\theta}(r)$ 

for  $r \in (0, \epsilon_x)$ ,  $\theta \in S^1$ 

♦ The metric :

$$\underbrace{\langle v, v \rangle = v_r^2 + f^2(r, \theta) v_{\theta}^2}_{\underbrace{}}$$

scalar product

length element

♦ Jacobi equation : second order linear ode

K determines M up to  $f_{rr} + Kf = 0$  $K(r, \theta)$  sectional curvature coverings

♦ Constant curvature :

 $f(r, \theta) = k^{-1} \sin(kr)$  $f(r, \theta) = k^{-1} \sinh(kr)$  $f(r, \theta) \sim r$  for small r positive curvature  $k^2$ negative curvature  $-k^2$ 

♦ Distance function (locally!!) :  $y = y_{\theta}(r) \implies d(x, y) = r$ 

♦ Hessian of distance function :

$$\nabla^2 r \simeq \begin{pmatrix} 0 & 0 \\ 0 & f_r/f \end{pmatrix} \qquad \frac{f_r/f \text{ may become}}{\text{negative or diverge !!}}$$

4/29

♦ Constant curvature :

"curvature of a sphere"  $S = f_r/f$ 

$$S(r, \theta) = k \cot(kr)$$

 $S(r, \theta) = k \coth(kr)$ 

positive curvature  $k^2$ 

negative curvature  $-k^2$ 

→ distance can fail to be convex or smooth!!

Comparison :

$$\alpha \leq K(r, \theta) \leq \beta \implies S_{\beta}(r) \leq S(r, \theta) \leq S_{\alpha}(r)$$

note reverse order

~> some local estimates

 $\frac{\pi}{2\sqrt{\beta}} \le r \le \frac{\pi}{2\sqrt{\alpha}} \qquad \frac{\pi}{\sqrt{\beta}} \le r \le \frac{\pi}{\sqrt{\alpha}}$ 

convexity is lost conjugate points (foci)

♦ Hessian of squared distance :

$$e_x(y) = d^2(x, y)$$
  $\nabla^2 e_x \simeq \begin{pmatrix} 2 & 0 \\ 0 & 2rS \end{pmatrix}$ 

♦ Constant curvature :

"curvature of a sphere"  $S = f_r/f$ 

$$S \sim r^{-1}$$
 for small  $r$   $S(r, \theta) = k \cot(kr)$ 

positive curvature  $k^2$ 

 $S(r, \theta) = k \coth(kr)$ 

negative curvature  $-k^2$ 

→ distance can fail to be convex or smooth!!

Comparison :

$$\alpha \leq K(r,\,\theta) \leq \beta \implies S_\beta(r) \leq S(r,\,\theta) \leq S_\alpha(r)$$

note reverse order

some local estimates

 $\frac{\pi}{2\sqrt{\beta}} \le r \le \frac{\pi}{2\sqrt{\alpha}} \qquad \frac{\pi}{\sqrt{\beta}} \le r \le \frac{\pi}{\sqrt{\alpha}}$ 

convexity is lost conjugate points (foci)

Hessian of squared distance :

$$e_x(y) = d^2(x, y)$$
  $\nabla^2 e_x \simeq \begin{pmatrix} 2 & 0 \\ 0 & 2rS \end{pmatrix}$ 

♦ Constant curvature :

"curvature of a sphere"  $S = f_r/f$ 

$$S \sim r^{-1}$$
 for small  $r$   $S(r, \theta) = k \cot(kr)$ 

positive curvature  $k^2$ 

negative curvature  $-k^2$ 

→ distance can fail to be convex or smooth!!

Comparison :

$$\alpha \leq K(r,\,\theta) \leq \beta \implies S_\beta(r) \leq S(r,\,\theta) \leq S_\alpha(r)$$

note reverse order

 $S(r, \theta) = k \coth(kr)$ 

some local estimates

non-positive curvature is well behaved

 $\frac{\pi}{2\sqrt{\beta}} \le r \le \frac{\pi}{2\sqrt{\alpha}} \qquad \frac{\pi}{\sqrt{\beta}} \le r \le \frac{\pi}{\sqrt{\alpha}}$ 

convexity is lost conjugate points (foci)

Hessian of squared distance :

$$e_x(y) = d^2(x, y)$$
  $\nabla^2 e_x \simeq \begin{pmatrix} 2 & 0 \\ 0 & 2rS \end{pmatrix}$ 

♦ Constant curvature :

"curvature of a sphere"  $S = f_r/f$ 

$$S \sim r^{-1}$$
 for small  $r$   $S(r, \theta) = k \cot(kr)$ 

positive curvature  $k^2$ 

 $S(r, \theta) = k \coth(kr)$ 

negative curvature  $-k^2$ 

→ distance can fail to be convex or smooth!!

Comparison :

$$\alpha \leq K(r,\,\theta) \leq \beta \implies S_\beta(r) \leq S(r,\,\theta) \leq S_\alpha(r)$$

note reverse order

some local estimates

non-positive curvature is well behaved

 $\frac{\pi}{2\sqrt{\beta}} \le r \le \frac{\pi}{2\sqrt{\alpha}} \qquad \frac{\pi}{\sqrt{\beta}} \le r \le \frac{\pi}{\sqrt{\alpha}}$ 

convexity is lost conjugate points (foci)

Hessian of squared distance :

$$\nabla^2 e_{\chi} \simeq \left(\begin{array}{cc} 2 & 0 \\ 0 & 2 \end{array}\right) + \left(\begin{array}{cc} 0 & 0 \\ 0 & 2(rS-1) \end{array}\right)$$

# Triangle comparison

 $\diamond$  Geodesic triangle  $\Delta = xyz$  :

$$(zy)^{2} = \underbrace{\left[(zx)^{2} - 2(zx)(xy)\cos \angle zxy + (xy)^{2}\right]}_{\leftarrow} + \mathcal{E}(\Delta)$$

planar triangle

Error estimate :

$$(xy)^{2}(DS_{\beta}(D) - 1) \leq \mathcal{E}(\Delta) \leq (xy)^{2}(DS_{\alpha}(D) - 1)$$

♦ Sign of K :

 $K \ge 0 \implies \alpha = 0 \implies \mathcal{E}(\Delta) \le 0 \text{ so } (zy) \le (zy)_{\text{plane}}$ 

 $K \le 0 \implies \beta = 0 \implies \mathcal{E}(\Delta) \ge 0 \text{ so } (zy) \ge (zy)_{\text{plane}}$ 

♦ Other comparisons :

metric, area, sum of angles, ...

$$f_{\beta}(r) \leq f(r, \theta) \leq f_{\alpha}(r)$$

f = 0 before  $f_{\alpha}$  and after  $f_{\beta}$ 

second-order Taylor of  $f(\gamma(t)) = d^2(z, \gamma(t))$ 





### Cut and conjugate locus

Cut(x) = { $\gamma_{\theta}(r_c) | \gamma_{\theta}$  not minimising after  $r_c$ } Conj(x) = { $r(\theta) | f(r, \theta) = 0$  for the first time}

- Why does a geodesic fail to minimise?
  - Contains conjugate points : not a local minimum
  - It is not unique (broken geodesics don't minimise)



$$y \in \operatorname{Cut}(x) \Leftrightarrow y \in \operatorname{Conj}(x) \text{ or } \underbrace{y = c(1) = \gamma(1)}_{}$$

typical case

Cut locus and topology

true in any dimension

Cut locus is a negligeable set

$$M = D_x \cup \operatorname{Cut}(x)$$
  $\operatorname{Cut}(x)$  deformation retract of  $M - \{x\}$ 

Injectivity radius

$$i(x) = \inf_{y \operatorname{Cut}(x)} d(x, y) = d(x, \operatorname{Cut}(x))$$

Klingenberg

$$i(x) \geq \min\left\{\frac{\pi}{\sqrt{\beta}}, \frac{\ell}{2}\right\}$$

### Cut and conjugate locus

Cut(x) = { $\gamma_{\theta}(r_c) | \gamma_{\theta}$  not minimising after  $r_c$ } Conj(x) = { $r(\theta) | f(r, \theta) = 0$  for the first time}

Why does a geodesic fail to minimise?

- Contains conjugate points : not a local minimum

- It is not unique (broken geodesics don't minimise)



$$y \in \operatorname{Cut}(x) \Leftrightarrow y \in \operatorname{Conj}(x) \text{ or } \underbrace{y = c(1) = \gamma(1)}_{}$$

typical case

Cut locus and topology

true in any dimension

Cut locus is a negligeable set

 $M = D_x \cup \operatorname{Cut}(x)$   $\operatorname{Cut}(x)$  deformation retract of  $M - \{x\}$ 

Injectivity radius

$$i(x) = \inf_{y \operatorname{Cut}(x)} d(x, y) = d(x, \operatorname{Cut}(x))$$

Klingenberg

$$i(x) \ge \min\left\{\frac{\pi}{\sqrt{\beta}}, \frac{\ell}{2}\right\}$$
  $\ell$  length of shortest loop through x

### Cut and conjugate locus

Cut(x) = { $\gamma_{\theta}(r_c) | \gamma_{\theta}$  not minimising after  $r_c$ } Conj(x) = { $r(\theta) | f(r, \theta) = 0$  for the first time}

- Why does a geodesic fail to minimise?
  - Contains conjugate points : not a local minimum
  - It is not unique (broken geodesics don't minimise)



$$y \in \operatorname{Cut}(x) \Leftrightarrow y \in \operatorname{Conj}(x) \text{ or } \underbrace{y = c(1) = \gamma(1)}_{}$$

typical case

Cut locus and topology

true in any dimension

Cut locus is a negligeable set

 $M = D_x \cup \operatorname{Cut}(x)$   $\operatorname{Cut}(x)$  deformation retract of  $M - \{x\}$ 

Injectivity radius

$$i(x) = \inf_{y \operatorname{Cut}(x)} d(x, y) = d(x, \operatorname{Cut}(x))$$

Klingenberg

$$i(x) \ge \min\left\{\frac{\pi}{\sqrt{\beta}}, \frac{\ell}{2}\right\}$$
  $\ell$  length of shortest loop through x

# Convexity of squared distance

♦ Recall the Hessian :

$$\nabla^2 e_{\mathbf{x}} \simeq \begin{pmatrix} 2 & 0 \\ 0 & 2rS \end{pmatrix} \ge 2 \,\beta r \cot(\beta r) \quad \text{(positive curvature)}$$

 $\geq 2$  (non-positive curvature)

♦ But is it convex ?

The problem comes from closed geodesics

Convexity radius : when is a geodesic ball convex?

$$B(x, R)$$
 is convex if  $R \le \frac{i(M)}{2}$ ;  $R \le \frac{\pi}{2\sqrt{\beta}}$ 

 Hadamard-von Mangoldt : squared distance is globally strongly convex on simply-connected surfaces of negative curvature





# Convexity of squared distance

Recall the Hessian :

$$\nabla^2 e_{\mathbf{x}} \simeq \begin{pmatrix} 2 & 0 \\ 0 & 2rS \end{pmatrix} \ge 2 \,\beta r \cot(\beta r) \quad \text{(positive curvature)}$$

 $\geq 2$  (non-positive curvature)

♦ But is it convex ?

The problem comes from closed geodesics

Convexity radius : when is a geodesic ball convex?

$$i(M) = \inf_{x} i(x)$$
  $B(x, R)$  is convex if  $R \le \frac{i(M)}{2}$ ;  $R \le \frac{\pi}{2\sqrt{\beta}}$ 

 Hadamard-von Mangoldt : squared distance is globally strongly convex on simply-connected surfaces of negative curvature





# Plan

Introduction with surfaces

### 2 Higher dimension

Convex stochastic optimisation

A Riemannian recursive estimation

### Proposed reading

Why do we need a connection ?

 $\nabla_v X$  = derivative of vector field X along direction v

Levi-Civita connection :

 $\nabla_{V}\langle X,Y\rangle = \langle \nabla_{V}X,Y\rangle + \langle X,\nabla_{V}Y\rangle$ 

(zero torsion)

(compatible with the metric)

$$\nabla_X Y - \nabla_Y X = [X, Y]$$

♦ Koszul formula :

elasticity tensor skew-symmetric

Why do we need a connection ?

 $\nabla_v X$  = derivative of vector field X along direction v

♦ Levi-Civita connection :

 $\nabla_{V}\langle X,Y\rangle = \langle \nabla_{V}X,Y\rangle + \langle X,\nabla_{V}Y\rangle$ 

(zero torsion)

B(t)

u(t)

$$\nabla_X Y - \nabla_Y X = [X, Y]$$







elasticity tensor

or skew-symmetric

(compatible with the metric)

 $p_t = v(t)$ 

 $q_t$ 

♦ Example 1 :

#### embedded submanifold

 $(\nabla_v X)(x) = \Pi_x(D_v X)$  projection of component-wise derivative

♦ Example 2 :

geodesic spherical coordinates

$$\nabla_{\partial_r} \partial_r = 0 \qquad \nabla_{\partial_r} \partial_\theta = S(r,\theta) \partial_\theta$$
$$\nabla_{\partial_\theta} \partial_r = S(r,\theta) \partial_\theta \qquad \nabla_{\partial_\theta} \partial_\theta = -(f^2)_r \partial_r$$

Interpretation :

$$\underbrace{\nabla_{\partial_r}\partial_r = 0}_{\text{geodesic equation}} \qquad \underbrace{\nabla_{\partial_\theta}\partial_r = S(r,\theta)\partial_\theta}_{\text{Hessian of distance}}$$

♦ Example 1 :

#### embedded submanifold

 $(\nabla_v X)(x) = \prod_x (D_v X)$  projection of component-wise derivative

♦ Example 2 :

geodesic spherical coordinates

$$\nabla_{\partial_r}\partial_r = 0 \qquad \qquad \nabla_{\partial_r}\partial_\theta = S(r,\theta)\partial_\theta$$

$$\nabla_{\partial_{\theta}}\partial_r = S(r,\theta)\partial_{\theta}$$

 $\nabla_{\partial_{\theta}} \partial_{\theta} = -(f^2)_r \partial_r$ 

♦ Interpretation :





### Some definitions

♦ Geodesic equation :

 $\nabla_{\dot{Y}} \dot{Y} = 0$  zero-acceleration parameterised curve  $\diamond$  Gradient  $\nabla f$ :  $\underbrace{v \cdot f}_{\text{directional derivative}} = \underbrace{df(x) \cdot v}_{df(x):T_x M \to M} = \underbrace{\langle \nabla f, v \rangle_x}_{\text{gradient}}$   $\diamond$  Hessian  $\nabla^2 f$ :  $\underbrace{\nabla^2 f \cdot v = \nabla_v \nabla f(x)}_{\text{self-adjoint endomorphism}} \qquad \underbrace{\nabla^2 f \cdot (u, v) = \langle \nabla^2 f \cdot u, v \rangle}_{\text{symmetric bilinear form}}$ 

♦ Second-order Taylor :

$$f(\gamma(1)) = f(\gamma(0)) + \left\langle \nabla f, \dot{\gamma} \right\rangle_{\gamma(0)} + \left. \frac{1}{2} \nabla^2 f \cdot (\dot{\gamma}, \dot{\gamma}) \right|_{t^*} \quad \text{where } t^* \in (0, 1)$$

# The exponential map

#### ♦ Definition from ODE theory :

$$\nabla_{\dot{\gamma}}\dot{\gamma} = 0$$
 i.e.  $\left(\gamma(0) = x \text{ and } \dot{\gamma}(0) = v = y^i \partial_i\right) \implies \operatorname{Exp}_x(v) = \gamma(1) = y$ 

the solution is unique for given initial conditions

#### Completeness and Hopf-Rinow :

 $\operatorname{Exp}_{x}(v)$  is defined for all  $v \iff$  any  $x, y \in M$  connected by a minimising geodesic

♦ Normal coordinates : is the relation  $y \mapsto y^i$  unique?

 $M = D_x \cup \operatorname{Cut}(x)$  Exp diffeomorphism of  $D_x$ 

- $\diamond$  What happens on Cut(x) :
- $|d \operatorname{Exp}_{x}(v)| = 0$  conjugate point
- Exp is not bijective

# The exponential map

#### ♦ Definition from ODE theory :

$$\nabla_{\dot{\gamma}}\dot{\gamma} = 0$$
 i.e.  $\left(\gamma(0) = x \text{ and } \dot{\gamma}(0) = v = y^i \partial_i\right) \implies \operatorname{Exp}_x(v) = \gamma(1) = y$ 

the solution is unique for given initial conditions

♦ Completeness and Hopf-Rinow :

 $\operatorname{Exp}_{x}(v)$  is defined for all  $v \iff$  any  $x, y \in M$  connected by a minimising geodesic

♦ Normal coordinates : is the relation  $y \mapsto y^i$  unique?

 $M = D_x \cup \operatorname{Cut}(x)$  Exp diffeo

Exp diffeomorphism of  $D_x$ 

- ♦ What happens on Cut(*x*) :
- $|d \operatorname{Exp}_{x}(v)| = 0$  conjugate point
- Exp is not bijective



# The exponential map

#### ♦ Definition from ODE theory :

$$\nabla_{\dot{\gamma}}\dot{\gamma} = 0$$
 i.e.  $\left(\gamma(0) = x \text{ and } \dot{\gamma}(0) = v = y^i \partial_i\right) \implies \operatorname{Exp}_x(v) = \gamma(1) = y$ 

the solution is unique for given initial conditions

#### ♦ Completeness and Hopf-Rinow :

 $\operatorname{Exp}_{x}(v)$  is defined for all  $v \iff$  any  $x, y \in M$  connected by a minimising geodesic

♦ Normal coordinates : is the relation  $y \mapsto y^i$  unique?

 $M = D_x \cup \operatorname{Cut}(x)$  Exp diffeomorphism of  $D_x$ 

- $\diamond$  What happens on Cut(x) :
- $|d \operatorname{Exp}_{x}(v)| = 0$  conjugate point
- Exp is not bijective



## Convex sets and functions

 $\diamond$  Convex  $A \subset M$  :

 $x, y \in A$ : length-minimising  $\gamma(0) = x, \gamma(1) = y$  and  $\gamma(t) \in A$ 

A ball may fail to be convex !!

Convexity radius : (small balls are always convex)

$$R_{cx}(M) \ge \min\left\{\frac{i(M)}{2}, \frac{\pi}{2\sqrt{\beta}}\right\}$$



$$f \circ \gamma : [0, 1] \to \mathbb{R}$$
 is convex



♦ Example :

$$e_x : B(x, R) \to \mathbb{R}$$
 where  $e_x(y) = d^2(x, y)$ 

 $\rightsquigarrow$  if  $R < R_{cx}$  this function is convex

Characterisation :

A is convex and  $\nabla^2 \varphi(y) \ge 0$  for  $y \in A$ 

# Hadamard manifolds

M simply connected, complete, with sectional curvature  $\leq 0$ 

- ♦ Examples :
  - Euclidean space
  - Poincaré half plane
  - Cones of covariance matrices
- ♦ Nice properties :  $i(M) = \infty \rightarrow no$  conjugate points, no closed geodesics
  - ◊ Exp is a diffeomorphism
  - squared distance is smooth
  - squared distance is strongly convex

 $\nabla^2 e_x(y) \ge 1$  all x and y in M  $(e_x(y) = d^2(x, y))$ 

$$F(y) = \frac{1}{N} \sum_{n=1}^{N} d^2(x_n, y) \quad \rightsquigarrow \text{ unique minimum and stationary point } \hat{x}_N$$

smooth strongly convex function

# Plan

Introduction with surfaces

### 2 Higher dimension

Convex stochastic optimisation

A Riemannian recursive estimation

### Proposed reading

# Stochastic optimisation

♦ Loss function  $(L : M \to \mathbb{R})$  :

$$L(x) = \mathbb{E}_z \ell(x, z)$$
 or  $L(x) = \frac{1}{N} \sum_{n=1}^N \ell(x, z_n)$ 

♦ Main issues :

loss function unknown ; evaluation too costly

Idea : learn and optimise at the same time!!

- generate or observe  $z_n$  where n = 1, 2, ...

- follow the gradient on average  $x_{n+1} = \operatorname{Exp}_{x_n}(-\gamma_{n+1}\nabla \ell(x_n, z_{n+1}))$ 

Deterministic vs stochastic :

	find stationary point	find local min	local rate of convergence
deterministic	YES	NO	geometric
stochastic	YES	YES	harmonic

♦ Limit set : connected component of  $\{\nabla L = 0\} \cup \{\infty\}$ 

(for more, recall the capture theorem)

# Local rate of convergence

 $x_{n+1} = \operatorname{Exp}_{x_n} \left( -\gamma_{n+1} \nabla \ell(x_n, z_{n+1}) \right)$ 

♦ Assumptions :

- $(x_n) \subset D$  compact convex set
- exactly one stationary point  $x^* \in D$
- L is *µ*-strongly convex in D
- controle of moments of noise

strong convexity 
$$(x, y \in D)$$
:  $L(y) - L(x) \ge \langle \nabla L(x), \operatorname{Exp}_{x}^{-1}(y) \rangle_{x} + \frac{\mu}{2} d^{2}(x, y)$ 

(a convex function is above all its tangents)

$$\rightsquigarrow$$
 strong attraction :  $-\frac{\mu}{2}d^2(x, x^*) \ge \langle \nabla L(x), \operatorname{Exp}_x^{-1}(x^*) \rangle_x$ 

(attraction to x\* is super-linear)

# Local rate of convergence

 $x_{n+1} = \operatorname{Exp}_{x_n} \left( -\gamma_{n+1} \nabla \ell(x_n, z_{n+1}) \right) \rightsquigarrow set \ u_{n+1} = \gamma_{n+1} \nabla \ell(x_n, z_{n+1})$ 

Assumptions :

- $(x_n) \subset D$  compact convex set
- exactly one stationary point  $x^* \in D$
- L is *µ*-strongly convex in D
- controle of moments of noise

triangle comparison :  $d^{2}(x_{n+1}, x^{*}) \leq d^{2}(x_{n}, x^{*}) + 2\langle u_{n+1}, \operatorname{Exp}_{x_{n}}^{-1}(x^{*}) \rangle + DS_{\alpha}(D) ||u_{n+1}||^{2}$ conditional expectation :  $\mathbb{E}_{n} d^{2}(x_{n+1}, x^{*}) \leq d^{2}(x_{n}, x^{*}) + 2\gamma_{n+1} \langle \nabla L(x_{n}), \operatorname{Exp}_{x_{n}}^{-1}(x^{*}) \rangle + C\gamma_{n+1}^{2}$ strong attraction :  $\mathbb{E}_{n} d^{2}(x_{n+1}, x^{*}) \leq (1 - \gamma_{n+1}\mu) d^{2}(x_{n}, x^{*}) + C\gamma_{n+1}^{2}$ take expectation :  $\mathbb{E} d^{2}(x_{n+1}, x^{*}) \leq (1 - \gamma_{n+1}\mu) \mathbb{E} d^{2}(x_{n}, x^{*}) + C\gamma_{n+1}^{2}$  $\diamond$  A first conclusion :

we must take  $\limsup \gamma_{n+1} \mu < 1$ 

# The problem of tuning

Vsual choice of step-size :

$$\gamma_n = \frac{A}{n^{\alpha} + B}$$
 where  $\alpha \in (0, 1]$   $\alpha \uparrow 1$  stops the algorithm faster

♦ Local rate of convergence :

$$\mathbb{E} d^2(x_n, x^*) \leq C \gamma_n^{\beta}$$
 where  $\beta \in (0, 1)$ 

♦ Optimal rate :

$$\mathbb{E} d^2(x_n, x^*) \le C\gamma_n \text{ requires } A > \frac{\alpha}{\mu}$$

(Please note these are only local rates!)

 $\label{eq:conclusion} \begin{array}{l} \circ \mbox{ Conclusion }: \ast \mbox{ we need to know } \mu \mbox{ (spend money)} \\ & \ast \mbox{ we need to guess } \mu \mbox{ (spend time)} \\ & \ast \mbox{ convergence can be arbitrarily bad} \\ & \ast \mbox{ anyway, a small } \mu \mbox{ is a bad case} \end{array}$ 

 $\diamond$  Can we get around knowing  $\mu$  ?

vi there exist some very nice tricks

## Averaged stochastic gradient

♦ Maintain a constant step-size :

 $x_{n+1} = \operatorname{Exp}_{x_n}(-\gamma \nabla \ell(x_n, z_{n+1}))$   $\gamma$  constant (or slowly decreasing)

♦ Does this converge ?

a stationary Markov process (the question is convergence in law, or ergodicity) ~> somehow, we need to stabilise it

Recursive Riemannian average: (generalise the Polyak average)

$$\hat{x}_{n+1} = \hat{x}_n \#_{\frac{1}{n+1}} x_{n+1}$$

geodesic weighted average

In a Euclidean space, this reduces to

$$\hat{x}_{n+1} = \frac{n}{n+1}\hat{x}_n + \frac{1}{n+1}x_{n+1}$$

→ SGD becomes the input of a Riemannian AR(1)

### A digression about barycentres

#### Riemannian barycentre (Fréchet mean) :

$$\bar{x}$$
 any global minimum of  $\mathcal{E}(x) = \int_{\mathcal{M}} d^2(x, y) P(dy)$ 

apparently, just stochastic optimisation

♦ In Euclidean space :

$$\bar{x} = \int_{\mathcal{M}} y P(dy)$$

unique global minimiser

♦ Law of large numbers :

$$\hat{x}_n = \frac{1}{n} \sum_{m=1}^n x_m \to \bar{x}$$
  $\hat{x}_{n+1} = \frac{n}{n+1} \hat{x}_n + \frac{1}{n+1} x_{n+1}$ 

♦ General Riemannian case :

 $\mathcal{E}(x)$  is non-differentiable, non-convex, and has multiple minima !!

- Conclsuion :

Open problem : barycentre of a Markov chain

# Plan

Introduction with surfaces

### 2 Higher dimension

Convex stochastic optimisation

A Riemannian recursive estimation

### Proposed reading

# The problem of recursive estimation

#### Stimation/learning problem :

 $\begin{array}{ll} \mbox{minimise a statistical divergence} & \theta^* = \mbox{argmin} D(\mathsf{P}_{\rm true} | \mathsf{P}_{\theta}) \\ & \theta \in \Theta \mbox{ (model space, a manifold)} \\ & \mbox{no reason to think } \mathsf{P}_{\rm true} = \mathsf{P}_{\theta^*} \end{array}$ 

 $\rightsquigarrow$  but we do not know  $P_{true}$  in the first place!!

Empirical estimation : (example of KL divergence)

$$\mathsf{D}(\mathsf{P}_{\mathrm{true}}|\mathsf{P}_{\theta}) = \int \log\left[\frac{\mathsf{p}_{\mathrm{true}}(x)}{\mathsf{p}_{\theta}(x)}\right] d\mathsf{P}_{\mathrm{true}}(x) \approx \frac{1}{N} \sum_{n=1}^{N} \log \mathsf{p}_{\mathrm{true}}(x_n) - \frac{1}{N} \sum_{n=1}^{N} \log \mathsf{p}_{\theta}(x_n)$$

first term does not depend on  $\theta$ 

#### Drawbacks

changes the original minimisation problem recomputes from scratch with new samples not suitable to very complicated models

#### Advantages

consistent, asymptotically efficient uses established optimisation methods

♦ Recursive estimation : we try to have the same advantages without the drawbacks

# The problem of recursive estimation

#### Stimation/learning problem :

minimise a statistical divergence

$$\begin{aligned} \theta^* &= \operatorname{argmin} \mathsf{D}(\mathsf{P}_{\mathsf{true}} | \mathsf{P}_{\theta}) \\ \theta &\in \Theta \text{ (model space, a manifold)} \\ \mathsf{no reason to think } \mathsf{P}_{\mathsf{true}} &= \mathsf{P}_{\theta^*} \end{aligned}$$

 $\rightsquigarrow$  but we do not know P<sub>true</sub> in the first place!!

Empirical estimation : (example of KL divergence)

$$D(P_{\text{true}}|P_{\theta}) = \int \log\left[\frac{p_{\text{true}}(x)}{p_{\theta}(x)}\right] dP_{\text{true}}(x) \approx \underbrace{\frac{1}{N} \sum_{n=1}^{N} \log p_{\text{true}}(x_n) - \frac{1}{N} \sum_{n=1}^{N} \log p_{\theta}(x_n)}_{\text{true}}$$

first term does not depend on  $\boldsymbol{\theta}$ 

 Recursive estimation : we try to have the same advantages without the drawbacks

$$\theta_{n+1} = \varphi(\theta_n, \ldots, x_{n+1}) \qquad \lim \theta_n = \theta^*$$

## The Fisher information metric

♦ First definition :

a metric adapted to the divergence  $D(P_{\theta}|P_{\theta+d\theta}) = \frac{1}{2} ||d\theta||^2 + \dots$ 

♦ Is this really a metric ?

case of Kullback-Leibler) 
$$\|d\theta\|^2 = -\sum_a \sum_b \mathbb{E}_{\theta} \left(\frac{\partial^2 \log p_{\theta}}{\partial \theta^a \partial \theta^b}\right) d\theta^a d\theta^b$$
  
Rao's discovery  $\underline{\|d\theta\|^2 = \|d\theta'\|^2}$ 

invariance by reparameterisation

Second definition (Chentsov's theorem) :

a formula is not a definition

there is (essentially) a unique metric on  $\Theta$  invariant by sufficient statistics

 $D(P_{\theta}|P_{\theta+d\theta}) = D(P_{\theta} \circ \varphi | P_{\theta+d\theta} \circ \varphi) \qquad (\varphi \text{ sufficient statistic})$ 

 $\rightsquigarrow$  many computations become automatic...  $\rightsquigarrow$  explains the appearance of affine-invariance

### **Recursive estimation**

♦ Gradient flow :

 $\dot{\theta} = -\nabla_{\theta} D(P_{true} | P_{\theta})$  Limit set (forward)

{∞} or {stationary point} or {stationary infinite set}
→ we cannot run this dynamical system !!

Stochastic approximation :

$$\theta_{n+1} = \operatorname{Exp}_{\theta_n} (\gamma_{n+1} u(\theta_n, x_{n+1}))$$

Limit set (a.s.) :

$$\left. \begin{array}{l} \sum \gamma_n = \infty \;,\; \sum \gamma_n^2 < \infty \\ \mathbb{E}_{\text{true}} \; u(\theta, \, x) = -\nabla \mathsf{D}(\theta) \end{array} \right\} \implies \text{ same as above }$$

 $\diamond$  Reflected algorithm : introduce "walls" to avoid going to  $\{\infty\}$ 

♦ Unstable points :

 $\varepsilon_n(\theta) = u(\theta, x) - \mathbb{E}_{\text{true}} u(\theta, x) \text{ (approximation noise)}$ isotropic noise  $\implies \mathbb{P}_{\text{true}}(\theta_n \rightarrow \text{unstable point}) = 0$ 

### **Recursive estimation**

♦ Gradient flow :

 $\dot{\theta} = -\nabla_{\theta} D(P_{true} | P_{\theta})$  Limit set (forward)

{∞} or {stationary point} or {stationary infinite set}
→ we cannot run this dynamical system !!

Stochastic approximation :

$$\theta_{n+1} = \operatorname{Exp}_{\theta_n}(\gamma_{n+1} u(\theta_n, x_{n+1}))$$
 or any  $C^2$  retraction

Limit set (a.s.) :

$$\left. \begin{array}{l} \sum \gamma_n = \infty \,, \ \sum \gamma_n^2 < \infty \\ \mathbb{E}_{\text{true }} u(\theta, \, x) = -\nabla \mathsf{D}(\theta) \end{array} \right\} \implies \text{ same as above }$$

 $\diamond$  Reflected algorithm : introduce "walls" to avoid going to  $\{\infty\}$ 

♦ Unstable points :

 $\varepsilon_n(\theta) = u(\theta, x) - \mathbb{E}_{\text{true}} u(\theta, x)$  (approximation noise)

isotropic noise  $\implies \mathbb{P}_{true}(\theta_n \rightarrow unstable point) = 0$ 

# Local rate of convergence

Sehavior at stable point :

 $\nabla D(\theta^*) = 0$ ;  $\nabla^2 D(\theta^*) > 0$  (least eigenvalue  $\lambda$ )

 $\lambda$  depends on the choice of metric

 $\diamond$  Strong attraction : if  $\lambda > \mu > 0$  there exists open  $\Theta^*$  at  $\theta^*$ 

$$-\mu \, d^2(\theta, \theta^*) \ge \langle \nabla D(\theta), \operatorname{Exp}_{\theta}^{-1}(\theta^*) \rangle_{\theta} \qquad \text{for all } \theta \in \Theta^*$$

♦ Best achievable rate :

$$\gamma_n = \frac{A}{n} \text{ and } A > \frac{1}{2\mu} \implies d^2(\theta, \theta^*) = O\left(n^{-1}\right)$$

 $\diamond$  Automatic tuning : (assume  $\theta^* = \theta_{true}$ )

information metric  $\rightsquigarrow \lambda = 1$ 

## Local rate of convergence

Sehavior at stable point :

 $\nabla D(\theta^*) = 0$ ;  $\nabla^2 D(\theta^*) > 0$  (least eigenvalue  $\lambda$ ) (Hessian as bilinear form)

 $\lambda$  depends on the choice of metric

 $\diamond$  Strong attraction : if  $\lambda > \mu > 0$  there exists open  $\Theta^*$  at  $\theta^*$ 

 $-\mu \ d^2(\theta,\theta^*) \geq \langle \nabla D(\theta), \operatorname{Exp}_{\theta}^{-1}(\theta^*) \rangle_{\theta} \qquad \text{for all } \theta \in \Theta^*$ 

♦ Best achievable rate :

$$\gamma_n = \frac{A}{n} \text{ and } A > \frac{1}{2\mu} \implies d^2(\theta, \theta^*) = O\left(n^{-1}\right)$$

 $\diamond$  Automatic tuning : (assume  $\theta^* = \theta_{true}$ )

information metric  $\rightsquigarrow \lambda = 1$ 

# Asymptotic normality

Normalised error :

$$\xi_n = \sqrt{\gamma_n} \operatorname{Exp}_{\theta^*}^{-1}(\theta_n)$$

♦ The CLT :

$$\xi_n \implies \mathcal{N}(0, \Sigma)$$
 (expressed in o.n.b.)

♦ Lyapunov equation :

$$\begin{aligned} H\Sigma + \Sigma H &= -A^2 \Sigma^* \qquad \gamma_n = \frac{A}{n} \\ H &= \frac{1}{2} \operatorname{Id} - \nabla^2 D(\theta^*) \\ \Sigma^* &= \mathbb{E}_{\operatorname{true}} \left( \varepsilon_n(\theta^*) \otimes \varepsilon_n(\theta^*) \right) \end{aligned}$$

What does this mean ??!

(asymptotic behavior) 
$$d\xi(t) = H\xi(t)dt + \Sigma^{-\frac{1}{2}} dW(t)$$

linear attraction + white noise

1

Asymptotic efficiency

information metric ~~>

$$\Sigma = \Sigma^* = (I(\theta^*))^{-1}$$
$$d(\theta_n, \theta^*) \Rightarrow \chi^2_{\dim \Theta}$$

useful for change detection

# Unknown information metric

- Examples of "difficult models" :
  - mixture models
  - neural networks
  - FIM known but complicated
- ◊ "partial FIM" :

retain diagonal part and try to find  $\lambda$  and  $\Sigma$ 

Automatic tuning : (averaged stochastic gradient under some suitable metric)

 $\theta_{n+1} = \exp_{\theta_n}(-\gamma u(\theta_n, x_{n+1}))$   $\gamma$  constant (or slowly decreasing)  $\hat{\theta}_{n+1} = \hat{\theta}_n \#_{\frac{1}{n+1}} \hat{\theta}_{n+1}$  geodesic average

 $\diamond$  this guarantees  $O(n^{-1})$  convergence rate and asymptotic efficiency  $\diamond$  Exp and # need to be manageable (chose a symmetric geometry ..)

# Plan

Introduction with surfaces

### 2 Higher dimension

Convex stochastic optimisation

A Riemannian recursive estimation

### Proposed reading

# Proposed reading

#### Manifolds and Riemannian geometry :

- J.M. Lee : Introduction to Topological manifolds
- J.M. Lee : Introduction to Smooth Manifolds
- J.M. Lee : Introduction to Riemannian manifolds

#### Information geometry :

S.I Amari : Methods of information geometry ...... : Learn from the state of the art !!

#### ♦ Recursive estimation :

Nevilson & Hasminskii : Stochastic approximation and recursive estimation Marie Duflo : Algorithmes Stochastiques + Random iterative models

#### Riemannian recursive estimation :

Bonnabel : Stochastic gradient descent on Riemannian manifolds

Tripuraneni & *al* : Averaging stochastic gradient descent on Riemannian manifolds